# Pruning from Adaptive Regularization

Lars Kai Hansen
Carl Edward Rasmussen
*CONNECT, Electronics Institute, Technical University of Denmark, B 349,
DK-2800 Lyngby, Denmark*

Inspired by the recent upsurge of interest in Bayesian methods we consider adaptive regularization. A generalization based scheme for adaptation of regularization parameters is introduced and compared to Bayesian regularization. We show that pruning arises naturally within both adaptive regularization schemes. As model example we have chosen the simplest possible: estimating the mean of a random variable with known variance. Marked similarities are found between the two methods in that they both involve a "noise limit," below which they regularize with infinite weight decay, i.e., they *prune*. However, pruning is not always beneficial. We show explicitly that both methods in some cases may increase the generalization error. This corresponds to situations where the underlying assumptions of the regularizer are poorly matched to the environment.

## 1 Introduction ────────────

We believe in Ockham's Razor: the generalization error of a model with estimated parameters is decreased by constraining capacity to a minimum needed for capturing the rule (see, e.g., Thodberg 1991; Solla 1992). However, this minimum may be hard to define for nonlinear noisy systems where the rule is ill-defined.

Pruning is a popular tool for reducing model capacity and pruning schemes have been successfully applied to layered neural networks (LeCun *et al.* 1990; Thodberg 1991; Svarer *et al.* 1993). While pruning is a discrete decision process, regularization introduces soft constraints such as weight decay (Moody 1991). A common feature of these techniques is the need for control parameters: stop criteria for pruning and weight decays for regularization. In Svarer *et al.* (1993) a statistical stop criterion was developed for pruning of networks for regression problems.

Recently, MacKay reviewed a Bayesian approach to adaptive regularization in the context of neural networks, demonstrating that the *evidence*-based method can improve the generalization properties and he compared it to cross-validation (MacKay 1992a,b). Cross-validation is known to be rather noisy; hence methods based on statistical arguments

are recommended (see, e.g., Akaike 1969; Moody 1991; Hansen 1993). In this presentation we define such an alternative approach for adaptive regularization, and we test it on a case of ultimate simplicity, namely that of estimating the mean of a gaussian variable of known variance. Detailed insight can be obtained for this case.

In the course of comparing the two schemes, we have discovered a new feature of adaptive regularization. We find that both approaches involve a "noise limit," below which they regularize with infinite weight decay, i.e., they *prune*. This finding unifies the pruning and regularization approaches to capacity control.

## 2 Specification of the Toy Problem

The problem is defined as follows: consider a student–teacher setup based on a *teacher* (with parameter $\tilde{w}$) providing $N$ examples of the form

$$y_m = \tilde{w} + \nu_m \qquad \nu_m \sim \mathcal{N}(0, \sigma^2) \qquad m = 1, \dots, N \tag{2.1}$$

where the normal noise contributions are independent and have zero mean and common known variance $\sigma^2$. Based on the training data set, $D = \{y_m \mid m = 1, \dots, N\}$, the task of the *student* (with parameter $w$) is to infer the mean: $\tilde{w}$. The measure of success for the student is the generalization error,

$$E_G(w) = \int d\nu P(\nu)(\tilde{w} + \nu - w)^2 = (\tilde{w} - w)^2 + \sigma^2 \tag{2.2}$$

This is the expected error on a random new test example for the specific student weight as estimated on the given training set. In evaluating an estimation scheme, our measure will be the *average generalization error* obtained by averaging over all possible training sets of the given size $N$.

In the next three sections we consider three such estimators: the usual "empirical mean" obtained from maximum likelihood, an estimator based on MacKay's maximum evidence, and a novel estimator based on explicit minimization of the expected generalization error. In Section 6 we compute the average generalization error of each of the schemes, as a function of the teacher parameter.

## 3 Maximum Likelihood Estimation

The likelihood of the student parameter associated with the training set $D$ is

$$P(D \mid w) = \prod_{m=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(y_m - w)^2\right] \tag{3.1}$$

Maximizing this with respect to $w$ produces the standard maximum likelihood estimator, the *empirical mean*

$$w_{\mathrm{ML}}(D) = \frac{1}{N} \sum_{m=1}^{N} y_m \qquad (3.2)$$

This estimator is *unbiased*, i.e., the average value, averaging over all possible training sets, is the true mean $\tilde{w}$.

For a specific training set, the generalization error is obtained using 2.2. However, as mentioned above, we are not so much interested in the value for the particular training set. Rather we are interested in the expected error, obtained by averaging over all possible training sets of size $N$. To compute this quantity, we note that $w_{\mathrm{ML}} \sim \mathcal{N}(\tilde{w}, \sigma^2/N)$, hence the average generalization error is

$$\langle E_G(w_{\mathrm{ML}}) \rangle_N = \langle (\tilde{w} - w_{\mathrm{ML}})^2 \rangle + \sigma^2 = \frac{\sigma^2}{N} + \sigma^2 \qquad (3.3)$$

The first term is the *average excess error* made by the student due to the finite training set. The second term is due to the imperfectness of the data; there is noise in the test examples used for grading the student. The minimal error approached for large training sets is simply the noise level.

## 4 Bayesian Regularization

A Bayesian scheme, recently suggested to the neural net community by MacKay, adopts two levels of inference. The first level of inference consists of estimating the teacher parameters from the data, conditioned on a parameterized *prior* distribution of the teacher mean value. The second level of inference consists of estimating the parameters of the prior, with the purpose of maximizing the *evidence*.

The probability $P(w \mid D)$ of the parameter conditioned on the data can be obtained using Bayes' rule,

$$P(w \mid D) = \frac{P(D \mid w)P(w)}{P(D)} \qquad (4.1)$$

by specifying a prior distribution $P(w)$ of the parameter. We follow MacKay and employ a parameterized prior: $P(w) \equiv P(w \mid \alpha)$, with a parameter $\alpha$ playing the role of a "weight decay," to be determined at the second level of inference. The prior takes the form

$$P(w \mid \alpha) = \sqrt{\frac{\alpha}{2\pi}} \exp\left(-\frac{1}{2}\alpha w^2\right) \qquad (4.2)$$

Using the likehood from equation 3.1, we arrive at the *posterior* distribution

$$P(w \mid D, \alpha) = P(D \mid \alpha)^{-1} \sqrt{\frac{\alpha}{2\pi}} \exp\left(-\frac{1}{2}\alpha w^2\right)$$

$$\times \prod_{m=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(y_m - w)^2\right] \tag{4.3}$$

To find the most probable teacher parameter, we maximize $P(w \mid D, \alpha)$ with respect to $w$, to get

$$w_{\text{MP}}(D, \alpha) = \frac{N}{N + \alpha\sigma^2} w_{\text{ML}}(D) \tag{4.4}$$

In the following we will suppress the explicit dependence on the training data.

The second level of inference is to estimate the regularization parameter, and we do this by again invoking Bayes' rule

$$P(\alpha \mid D) = \frac{P(D \mid \alpha)P(\alpha)}{P(D)} \tag{4.5}$$

We assume the prior on $\alpha$ to be flat,[1] indicating that we do not know what value $\alpha$ should take. Hence, the most probable regularization is obtained by maximizing the likelihood of it. MacKay dubbed this quantity the *evidence*

$$P(D \mid \alpha) = \int dw P(D \mid w, \alpha)P(w \mid \alpha) \tag{4.6}$$

$$= \int_{-\infty}^{\infty} dw \sqrt{\frac{\alpha}{2\pi}} \exp\left(-\frac{1}{2}\alpha w^2\right)$$

$$\times \prod_{m=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(y_m - w)^2\right]$$

$$\Longrightarrow \log P(D \mid \alpha) = \frac{1}{2}\log(\alpha/2\pi) - \frac{1}{2}\log\left(\frac{N + \alpha\sigma^2}{2\sigma^2}\right)$$

$$+ \frac{(Nw_{\text{ML}})^2}{2\sigma^2(N + \alpha\sigma^2)} + \text{const} \tag{4.7}$$

where we have lumped all terms not depending on $\alpha$ in the constant. Maximizing the evidence with respect to $\alpha$,

$$\alpha_{\text{ML}} = \begin{cases} \dfrac{1}{w_{\text{ML}}^2 - \sigma^2/N} & w_{\text{ML}}^2 > \sigma^2/N \\ \infty & w_{\text{ML}}^2 \leq \sigma^2/N \end{cases} \tag{4.8}$$

---

[1] By flat we mean flat over $\log(\alpha)$, since $\alpha$ is a scale parameter.

which may be substituted into the expression for $w_{\mathrm{MP}}$ to find

$$
w_{\mathrm{MP}} = \begin{cases} \left(1 - \dfrac{\sigma^2}{N w_{\mathrm{ML}}^2}\right) w_{\mathrm{ML}} & w_{\mathrm{ML}}^2 > \sigma^2/N \\ 0 & w_{\mathrm{ML}}^2 \le \sigma^2/N \end{cases} \tag{4.9}
$$

Hence there is a sharp "noise limit" below which the scheme tells us to "prune," that is, to regularize with infinite strength. The noise limit has a straightforward interpretation: prune the parameter if the signal-to-noise ratio is less than unity. Note that the estimator is biased; its mean value is not $\tilde{w}$. It is *consistent* since $w_{\mathrm{MP}} \to \tilde{w}$ for large samples, while the noise limit shrinks to zero. To illustrate the power of the method, we compute in Section 6 the (training set) ensemble average of the generalization error.

In the above derivation we have used the "flat in $\log(\alpha)$ space" prior. This prior is improper, since it cannot be normalized. However, this is not important for the evidence approximation. To see this we introduce limits: $\alpha \in [\alpha_1, \alpha_2]$. The value of $w_{\mathrm{MP}}$ will not be affected by these limits as long as $\alpha_{\mathrm{ML}} \in [\alpha_1, \alpha_2]$. Choosing $\alpha_1 \ll N/\sigma^2$ and $\alpha_2 \gg N/\sigma^2$ the effect of the limits will be negligible for all $\alpha$. This is seen by introducing the limits in equation 4.4. In particular, we note that qualitative pruning still takes place: $|w_{\mathrm{MP}}| = |w_{\mathrm{ML}}| N/(N + \alpha_2 \sigma^2) \ll |w_{\mathrm{ML}}|$ for $w_{\mathrm{ML}}^2 < \sigma^2/N + 1/\alpha_2$.

The evidence framework is an approximation to the ideal Bayesian approach. Instead of using the maximum likelihood value of $\alpha$, we should ideally integrate over the distribution of $\alpha$, i.e., evaluate $P(w) = \int P(w \mid \alpha) P(\alpha) d\alpha$ and then use the *posterior mean* for our prediction: $w_{\mathrm{PM}} = \int w P(w) dw$. In contrast to the evidence framework the ideal approach is quite sensitive to the upper limit $\alpha_2$. Indeed if $\alpha_2 = \infty$ then $w_{\mathrm{PM}} = 0$ regardless of the value of $w_{\mathrm{ML}}$. For a range of $\alpha_2$ we find that $w_{\mathrm{PM}}$ shows a qualitative pruning behavior similar to that of $w_{\mathrm{MP}}$.

## 5 Generalization Error Minimization

While the evidence is an interesting statistical quantity, it is not obvious what the relation is between maximizing the evidence and minimizing test error (MacKay 1992a). Since the latter often is the basic modeling objective, we propose here to base the optimization of $\alpha$ on the expected generalization error (cf. equation 2.2). A similar approach was mentioned in Moody (1991).

In order to be able to compare directly with the Bayesian approach, we use a regularized *least-squares* procedure to find the student weight

$$
E_{\mathrm{T}} = \frac{1}{2\sigma^2} \sum_{m=1}^{N} (y_m - w)^2 + \frac{1}{2} \alpha w^2 \tag{5.1}
$$

Minimizing with respect to $w$, we recover 4.4

$$
w_{\mathrm{G}} = \frac{N}{N + \alpha \sigma^2} w_{\mathrm{ML}} \tag{5.2}
$$

Our aim is to minimize the average generalization error (averaged again with respect to training sets). To this end we note that the distribution of $w_G(D, \alpha)$ can be computed using simple manipulations of random variables. Noting again that $w_{ML} \sim \mathcal{N}(\tilde{w}, \sigma^2/N)$ we have $w_G \sim \mathcal{N}\left(\frac{N}{N+\alpha\sigma^2}\tilde{w}, \frac{N}{(N+\alpha\sigma^2)^2}\sigma^2\right)$.[2] Consequently the averaged generalization error becomes

$$\langle E_G(w_G) \rangle_N = \left(\frac{\alpha\sigma^2}{N+\alpha\sigma^2}\tilde{w}\right)^2 + \left[1 + \frac{N}{(N+\alpha\sigma^2)^2}\right]\sigma^2 \tag{5.3}$$

You might be uncomfortable with the appearance of the unknown teacher parameter $\tilde{w}$, however, it will shortly be replaced by an estimate based on the data. Minimizing the generalization error 5.3 with respect to $\alpha$, we find the simple expression

$$\alpha = \frac{1}{\tilde{w}^2} \tag{5.4}$$

which is inserted in equation 5.2, to obtain

$$w_G = \frac{\tilde{w}^2}{\tilde{w}^2 + \sigma^2/N} w_{ML} \tag{5.5}$$

To proceed we need to insert an estimate of the teacher parameter $\tilde{w}$. We will consider two cases; first we try setting $\tilde{w} = w_{ML}$ in equation 5.5, hence

$$w_{G1} = \frac{w_{ML}^2}{w_{ML}^2 + \sigma^2/N} w_{ML} \tag{5.6}$$

This estimator is biased, but consistent and it does not call for pruning even if the data are very noisy.

Secondly, being more brave, we could argue that the best estimate we have of $\tilde{w}$ is $w_G$ and accordingly set $\tilde{w} = w_{G2}$ in equation 5.5, to obtain a self-consistent equation[3] for $w_{G2}$

$$w_{G2} = \frac{w_{G2}^2}{w_{G2}^2 + \sigma^2/N} w_{ML} \tag{5.7}$$

By substituting $\tilde{w} = w_{G2}$, the student is not operating with the true cost function 5.3, but with a modified, self-consistent, cost function depicted in Figure 1. We note that by this substitution, a potentially dangerous global minimum is created at $w_{G2} = 0$.

---

[2]Using $(\xi = a + b\psi) \wedge [\psi \sim \mathcal{N}(c, d^2)] \Rightarrow \xi \sim \mathcal{N}(a + bc, b^2d^2)$.

[3]Interestingly, this equation is recovered from the Bayesian approach if we optimize $w$ and $\alpha$ simultaneously (i.e., if we do not separate inference in two distinct levels). In this case Bayes' rule becomes

$$P(w, \alpha \mid D) = \frac{P(D \mid w, \alpha)P(w \mid \alpha)P(\alpha)}{P(D)}$$
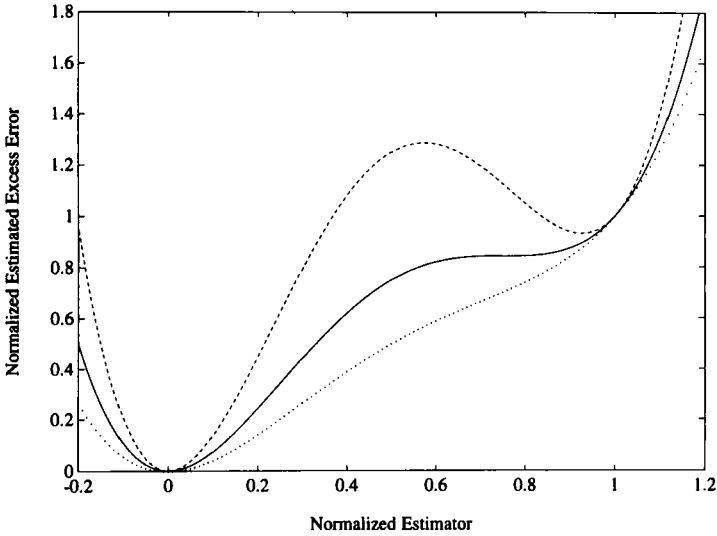
Figure 1: The estimated average excess generalization error as function of normalized estimator. The graphs are obtained by introducing equation 5.2 in 5.3, and setting $\tilde{w} = w_G$. The estimator is scaled by $w_{ML}$, and the estimated excess error is scaled by $\sigma^2/N$. The three graphs correspond to different noise levels (dashed, low; solid, critical; dotted, high). Note that if the noise is low, a local minimum is found, otherwise pruning takes place.

We could envision equation 5.7 solved by iteration, e.g., starting from $w_{G2}^{(0)} = w_{ML}$. Iterating the function on the right-hand side we see that, besides the solution $w_{G2} = 0$, there may be two more fixed points (or solutions) depending on the parameters. Analyzing the iteration scheme for this simple case is straightforward, one of the possible fixed points is stable, the other unstable. If it exists, the stable fixed point (corresponding to the local minimum in Fig. 1) is found by the iteration scheme, and it is given by

$$
w_{G2} = \begin{cases} \dfrac{w_{ML}}{2}\left(1 + \sqrt{1 - 4\sigma^2 \big/ N w_{ML}^2}\right) & w_{ML}^2 > 4\sigma^2/N \\ 0 & w_{ML}^2 \leq 4\sigma^2/N \end{cases} \tag{5.8}
$$

Note that the noise limits are different for the generalization error method and Bayes' method. The generalization error method is more conserva-
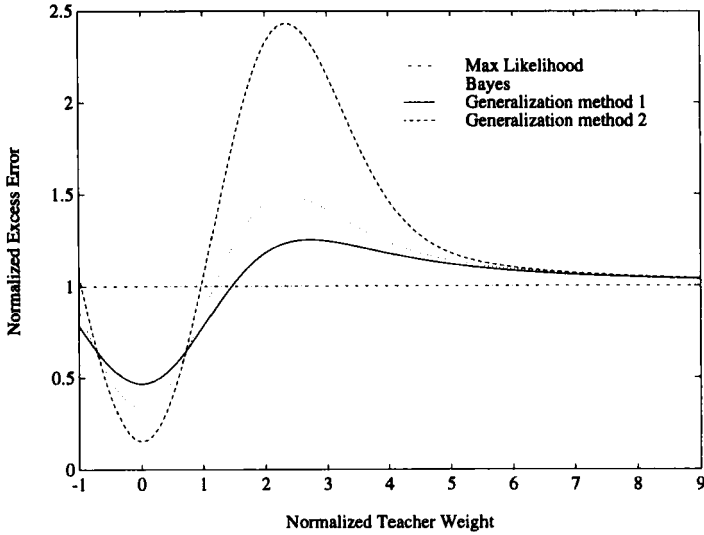
Figure 2: Average excess generalization error $\Delta E = E - \sigma^2$ as a function of teacher magnitude for the two adaptive regularization schemes used to estimate the mean of a gaussian variable. The teacher weight has been scaled by $\sqrt{\sigma^2/N}$, and the average excess error is scaled by the excess error of the maximum likelihood estimator ($\sigma^2/N$).

tive than Bayes' (if we associate conservatism with setting $w = 0$). The noise limit is a factor 4 larger, meaning that pruning or super-regularization will happen more often. Note also that this estimator, like the Bayesian, is biased but consistent.

## 6 Comparing the Regularization Schemes

As discussed above, the quantity to compare is the average (with respect to training sets) generalization error of the estimators. Since each estimator is a function of the empirical mean, $w_{ML}$, only, the average involves averaging with respect to the distribution of this quantity and we already noted that $w_{ML} \sim \mathcal{N}(\tilde{w}, \sigma^2/N)$. The generalization error of the maximum likelihood estimator itself was given in 3.3; it is independent of the teacher magnitude.

Since the regularized estimators depend only on the training set value of $w_{ML}$, the average excess generalization errors $[\Delta E(\tilde{w}) = E(\tilde{w}) - \sigma^2]$ are functions of $\tilde{w}/\sqrt{\sigma^2/N}$

$$\Delta E(\tilde{w}) = \frac{1}{\sqrt{2\pi\sigma^2/N}} \int_{-\infty}^{\infty} dw_{ML} \exp\left[-\frac{(w_{ML} - \tilde{w})^2}{2\sigma^2/N}\right] [\tilde{w} - w_E(w_{ML})]^2 \quad (6.1)$$

where $w_E$ is either of the three biased estimators. The integral is easily evaluated numerically. In Figure 2 we picture these functions, and we note that adaptation of the regularization parameter leads to a decreased error *only if the teacher is small*. If the teacher is very large the adaptive schemes lead to a negligible increase in error as they indeed should. However, if the teacher is of *intermediate* size all schemes produce seriously increased errors. The reason for this increased error is simply that the adaptive schemes are too conservative; there is a certain probability even for a unit magnitude teacher that the scheme will prune, hence introduce a large error.

We conclude that the benefit of the adaptive schemes depends critically on the distribution of teachers, i.e., the extent to which the domain lives up to the assumptions of the regularizer. In fact, distributions of teachers can be given for which each of the estimation schemes would be preferred. We note that for a teacher distribution that has significant mass at small teachers like the $P_0(\tilde{w}) \sim 1/|\tilde{w}|$ distribution corresponding to the prior of the Bayesian scheme (Buntine and Weigend 1991) the *brave generalization student* is superior for the present problem. This fact is most easily appreciated by noting that integrating the functions in Figure 2 with respect to $P_0(\tilde{w})$ corresponds to integrating the functions on a logarithmic ordinate axis with a uniform measure, hence stretching the region around $|\tilde{w}| \sim 0$ to infinity. We iterate, this is no big surprise since that student precisely minimizes the generalization error with the correct (implicit) prior.

Our experience with adaptive regularization is generally positive. A detailed account of multivariate applications, where pruning is achieved by utilizing individual weight regularization, will be given elsewhere (Rasmussen and Hansen 1994). The idea of using the *estimated* test error as a cost function for determination of regularization parameters is quite general. We are currently pursuing this in the context of time series processing with feedforward networks for which the corresponding generalization error estimate is well-known (see, e.g., Svarer *et al.* 1993).

In conclusion, we have shown that the use of adaptive regularization for the simple task of estimating the mean of a gaussian variable of known variance leads to nontrivial decision processes. We have defined a competitive scheme for identification of the regularization parameter, based on minimalization of the expected generalization error. We have shown that pruning can be viewed as infinite regularization, and that this is a natural consequence of adaptive regularization.

## Acknowledgments

## References

Akaike, H. 1969. Fitting autoregressive models for prediction. *Ann. Inst. Stat. Mat.* **21**, 243–247.

Buntine, W., and Weigend, A. 1991. Bayesian back-propagation. *Complex Syst.* **5**, 603–643.

Hansen, L. 1993. Stochastic linear learning: Exact test and training error averages. *Neural Networks* **6**, 393–396.

Le Cun, Y., Denker, J., and Solla, S. 1990. Optimal brain damage. In *Advances in Neural Information Processing Systems*, D. Touretzky, ed., Vol. 2, pp. 598–605. Morgan Kaufmann, San Mateo, CA.

MacKay, D. 1992a. Bayesian interpolation. *Neural Comp.* **4**, 415–447.

MacKay, D. 1992b. A practical framework for backpropagation networks. *Neural Comp.* **4**, 448–472.

Moody, J. 1991. Note on generalization, regularization and architecture selection in nonlinear systems. In *Neural Networks for Signal Processing I*, S. Kung, B. Juang, and C. Kamm, eds., pp. 1–10. IEEE, Piscataway, NJ.

Rasmussen, C., and Hansen, L. 1994. In preparation.

Solla, S. 1992. Capacity control in classifiers for pattern recognizers. In *Neural Networks for Signal Processing II*, S. Y. Kung *et al.*, eds., pp. 255–266. IEEE, Piscataway, NJ.

Svarer, C., Hansen, L., and Larsen, J. 1993. On design and evaluation of tapped-delay neural network architectures. In *IEEE International Conference on Neural Networks*, H. R. Berenji *et al.* eds., pp. 46–51. IEEE, Piscataway, NJ.

Thodberg, H. 1991. Improving generalization on neural networks through pruning. *Int. J. Neural Syst.* **1**, 317–326.